

Submission in Response to NSF CI 2030 Request for Information

DATE AND TIME: 2017-04-05 16:30:36

PAGE 1

REFERENCE NO: 288

This contribution was submitted to the National Science Foundation as part of the NSF CI 2030 planning activity through an NSF Request for Information, https://www.nsf.gov/publications/pub_summ.jsp?ods_key=nsf17031. Consideration of this contribution in NSF's planning process and any NSF-provided public accessibility of this document does not constitute approval of the content by NSF or the US Government. The opinions and views expressed herein are those of the author(s) and do not necessarily reflect those of the NSF or the US Government. The content of this submission is protected by the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License (<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>).

Author Names & Affiliations

- Brian Glendenning - NRAO

Contact Email Address (for NSF use only)

(Hidden)

Research Domain, discipline, and sub-discipline

Math and Physical Sciences, Astronomy, Radio Astronomy

Title of Submission

National Radio Astronomy Observatory Input to NSF CI 2030

Abstract (maximum ~200 words).

The National Radio Astronomy Observatory operates the Karl G. Jansky Very Large Array (VLA) and is the Executive of the North American part of the Atacama Large Millimeter/Submillimeter array. Hundreds of PI groups per year get observing time on these telescopes, and our accumulated historical data is frequently used for Archival research. We suggest that NSF investment in the following areas would be very helpful.

1. High-speed networking to remote (telescope) locations.
2. Long-duration, high-performance scientific software infrastructure.
3. Internationally interoperable portals, data flow, and computing resource allocation.
4. HPC Center focus, storage and duration.
5. Visualization and information extraction from multi-peta-pixel multi-dimensional image data.

Question 1 Research Challenge(s) (maximum ~1200 words): Describe current or emerging science or engineering research challenge(s), providing context in terms of recent research activities and standing questions in the field.

The National Radio Astronomy Observatory (NRAO) operates the Karl G. Jansky Very Large Array (VLA) near Socorro New Mexico, and is the operating partner (Executive) for the North American part of the Atacama Large Millimeter/Submillimeter Array (ALMA), which operates at a high site near San Pedro, Chile.

Both telescopes are very general purpose. Telescope time is allocated based on a peer-review process from many sub-fields of astronomy. Hundreds of PI groups per year get data, and in addition once the proprietary period has expired (usually one year), the data may be used

Submission in Response to NSF CI 2030 Request for Information

DATE AND TIME: 2017-04-05 16:30:36

PAGE 2

REFERENCE NO: 288

by other groups for Archival research.

Both telescopes are radio interferometers, which operate by coherently combining the signals of the relocatable antennas (27 for the VLA, 66 for ALMA) in complex central electronics (notably the correlators, which are approximately 0.1 Exa-Op very parallel special purpose supercomputers) which produces raw data, essentially a noisy (electronics, radio-frequency-interference, atmospheric and other environmental effects) irregularly sampled spatial Fourier transform of sky "stacked" over separate frequency channels for up to 4 polarizations.

The electronics are capable of sustaining 1 (VLA) and 16 (ALMA) Gigabytes per second of raw data output, although the data rates are usually averaged down (in time, and frequency) to a small fraction of that (of order 25 Megabytes/second for the VLA, and 6 for ALMA). This averaging is done both to reduce the computing that is needed, and because many times the science application do not need high data rates. However some do and there are some classes of science observations that are not made because computing capacity is not available.

The raw data is turned into regularly gridded 2-4 dimensional images (axes: position on the sky, frequency or Doppler velocity, polarization) using multi-million line of code software systems produced by the NRAO and our partners. These images (currently: Giga-pixel, coming Tera-pixel, Possible: Peta-pixel) are then typically processed through analysis codes (both produced by NRAO and the wider community) to enable the science to be extracted from the data.

Question 2 Cyberinfrastructure Needed to Address the Research Challenge(s) (maximum ~1200 words): Describe any limitations or absence of existing cyberinfrastructure, and/or specific technical advancements in cyberinfrastructure (e.g. advanced computing, data infrastructure, software infrastructure, applications, networking, cybersecurity), that must be addressed to accomplish the identified research challenge(s).

1. High-speed networking to remote (telescope) locations. To enable high data rate observing, 100 Gb class links to remote telescope sites need to be available. (There is buffer storage at the telescope locations, but the main processing centers are elsewhere).
2. Long-duration, high-performance scientific software infrastructure. Our domain specific software is a capital expense, and the hardware it runs on is an operational expense. Our software is built up over decades (we have two large systems in production at the moment, a Fortran based system approaching 40 years old, and a C++ based system approaching 25) with hundreds of FTE-years of effort invested, some of it from individuals with very particular expertise in the underlying calibration and imaging algorithms. The underlying computing paradigm changes many times during the lifetime of these software systems, and consequently it can be difficult for us to maintain good performance characteristics over time and over a wide range of platforms; our software runs on laptops, servers, clusters, commercial cloud providers, and XSEDE (currently minor usage, intend to extend). It would be extremely helpful if there was a scientific software middleware that computes on complex scientific data structures (not just matrices) in parallel, including I/O, with continued porting support as the computing landscape changes. This would allow us to concentrate on radio astronomy algorithms while retaining high performance as the computing landscape changes from underneath us.
3. Internationally interoperable portals, data flow, and computing resource allocation. Astronomy is an international activity. For example, the US is a minority partner of the ALMA Observatory. Many of the PI teams are international. At the moment national computing infrastructures tend to be stovepiped, making it difficult for us to support our users, who must navigate issues related to identity, data flow between systems, and computing resource allocation. International standards, agreements, and de-facto real-world international interoperability of these aspects to reduce the friction for our users would be very helpful.
4. HPC Center focus, storage and duration. Our use case is much more high throughput than high-performance computing. Our algorithms have relatively low FLOPS/IO ratios, so we need relatively large local high-speed scratch storage compared to typical HPC center provisioning. Different aspects of the processing have different needs (memory, I/O, FLOPS, parallelization), so having access to heterogeneous systems so we can efficiently process each stage without having to make an internet data movement would be very helpful. Some interactive access so that the data processing results can be evaluated before initiating the next processing stage is often required. As for permanent (archival) storage, if NSF cyberinfrastructure is to play a role it must be possible to make durable agreements over decades. We have all of our data from the beginning of VLA science operations in the late 1970s, and for us to move it out-of-house would require durable agreements. This is not just a matter of saving all the bits (old data is small since it was generated with Moore's law deflated detectors and computing systems), but also that it can be understood in modern software systems, i.e. the semantics of the data. Long term data accessibility matters involve file formats, versioning, documentation, standards, and knowledge embedded in software systems.

Submission in Response to NSF CI 2030 Request for Information

DATE AND TIME: 2017-04-05 16:30:36

PAGE 3

REFERENCE NO: 288

5. Visualization and information extraction from multi-peta-pixel multi-dimensional image data. Our images are currently typically in the Gigapixel range, with Terapixel images unusual but becoming more common, and Petapixel class images starting to be considered. (The forthcoming (starting September 2017) VLA Sky Survey would generate multi-Petapixel images if we had the computing and permanent storage to support it). We do not have good tools to visualize or extract information from the larger images. Our images are unusual in that the third (frequency/Doppler) axis is often deep, 10k or more pixels, and we sometimes have a shallow 4th axis (1-4 polarizations). The radio-astronomical community does not have good tools to handle images of this size. We have not yet found a tool from another community (e.g., HEP, GIS) that is suitable for our data. (Our in-house software is suitable for smaller, i.e. Gigapixel, images).

Question 3 Other considerations (maximum ~1200 words, optional): Any other relevant aspects, such as organization, process, learning and workforce development, access, and sustainability, that need to be addressed; or any other issues that NSF should consider.

As an NSF Large Facility (FFRDC), it is clear that NRAO has similar concerns to peer institutes over workforce development and talent management. The need to attract and retain world-class resources in direct competition with Universities, National Labs and Industry, particularly for software, HPC, and engineering positions, presents a key challenge. We see good value in the workshops and facilities meetings that NSF sponsors to understand these issues, but remain vulnerable to turnover and loss of highly specialized skills.

Consent Statement

- "I hereby agree to give the National Science Foundation (NSF) the right to use this information for the purposes stated above and to display it on a publically available website, consistent with the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License (<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>)."
-